# A Distributed Approach to Data Sharing:

## *DSSTox Public Toxicity Database Network*

Ann Richard

Environmental Carcinogenesis Division
National Health & Environmental Effects Research Lab
US Environmental Protection Agency

# Agency Problem:

- ➢ Too many chemicals to test
- ➢ Lack of sufficient and relevant data

Prioritize and focus limited resources on chemicals and problem areas estimated to pose greatest potential hazard

Inerts

TSCA/PMN

Endocrine Disruption
Testing Program

HPV Testing Pgm

Water CCL

Pesticides

# Computational Toxicology:

The application of mathematical and computer models and molecular biological approaches to improve EPA's

➢ prioritization of data requirements

➢ risk assessments

Chemicals of concern → Gather relevant information → Develop and apply models → Perform extrapolations → Inform risk assessment

# Chemistry-based Data Mining & Exploration:

**Chemical(s) of concern**

*Chemical-specific data*

*Structural analogs*

*Property analogs*

*Biological/ mechanistic analogs*

**Structure-Activity Relationships**

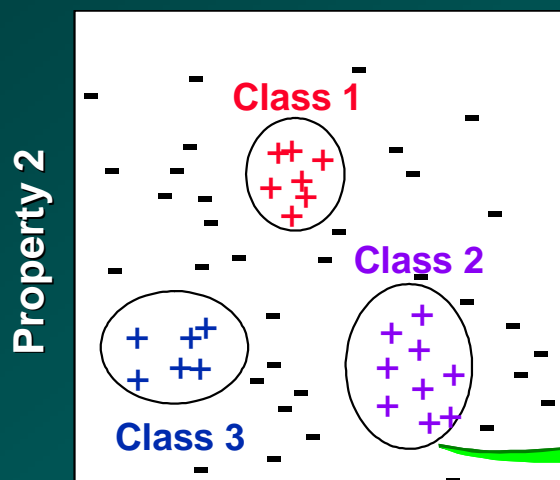# Structure-Activity Relationships (SAR):

## Activity = f (structure)

*Analogy*
*Heuristics*
*Machine-Learning Inference*
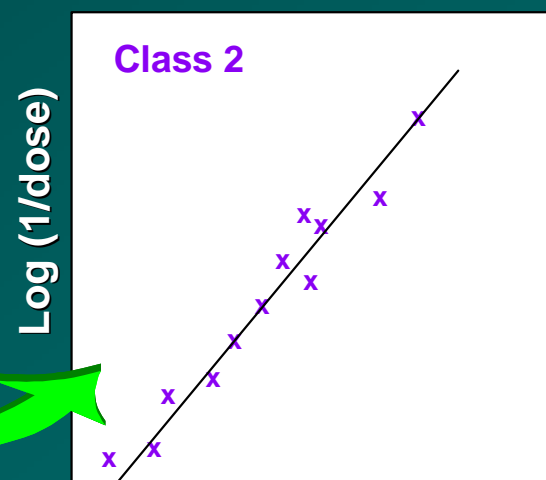*Statistical Correlation:*



SAR Classification
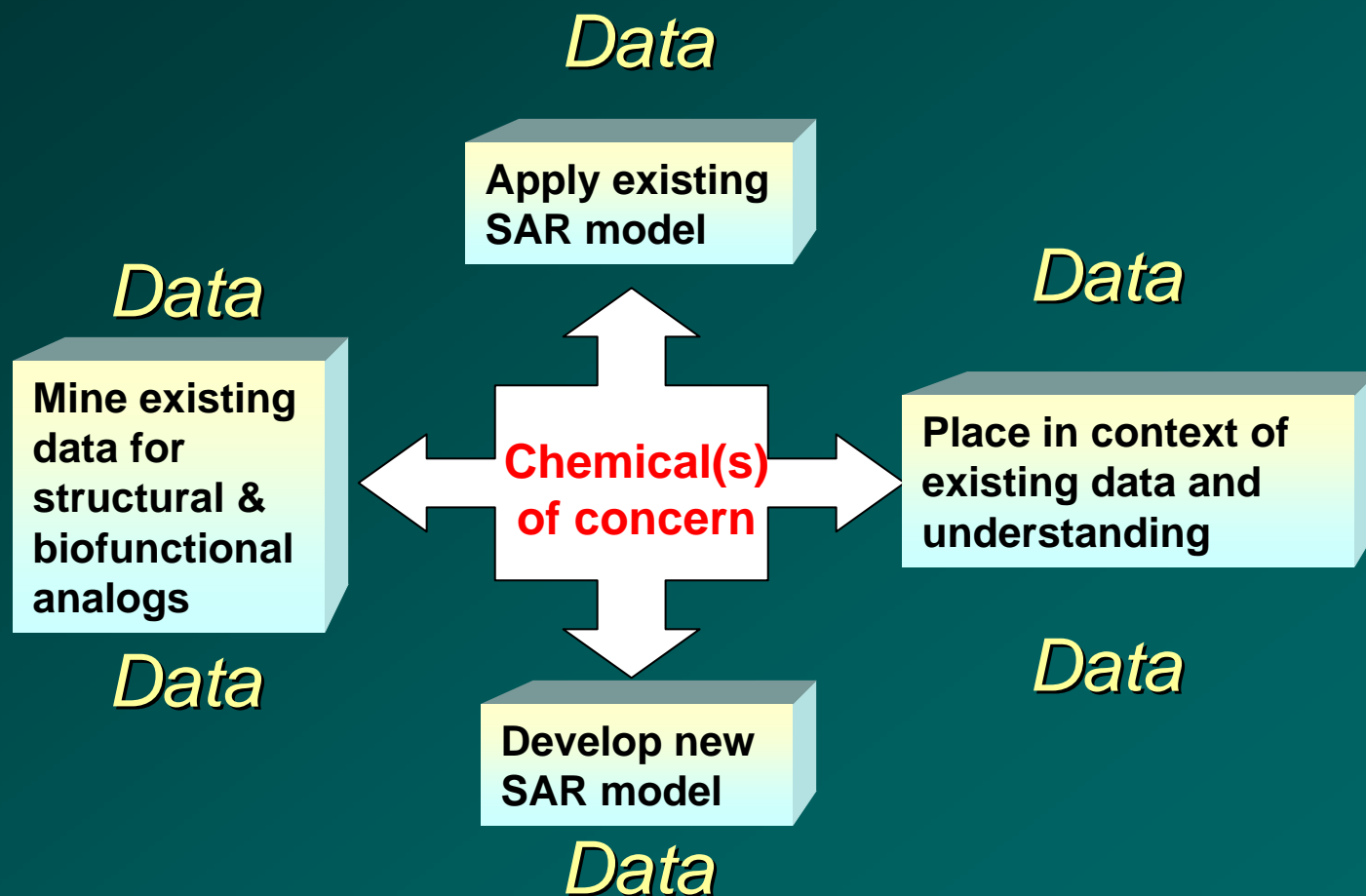
QSAR Correlation

# Structure-based Screening & Prioritization:

# Limitations of Public Toxicity Data for Use in SAR:

- Scattered sources

- Non-standard formats

- Diverse information content

- Lack of chemical structure annotation

- Cannot access full database

Chemical Structures

+

Public Toxicity Data Files

**D**istributed
**S**tructure-**S**earchable
**To**xicity
Public
Database
Network

DSSTox SDF Files

DSSTox Users

Import into User Database Applications

Structural Analog Searching

Improved Toxicity Prediction Models

## U.S. Environmental Protection Agency

# Distributed Structure-Searchable Toxicity (DSSTox) Public Database Network

Recent Additions | Contact Us | Print Version        Search: [        ] **GO**

EPA Home > Research & Development > National Health and Environmental Effects Research Laboratory > Distributed Structure-Searchable Toxicity (DSSTox) Public Database Network

**About DSSTox**

**Work in Progress**

**Frequent Questions**

**Databases**

**Central Field Definition Table**

**Apps, Tools & More**

**DSSTox Community**
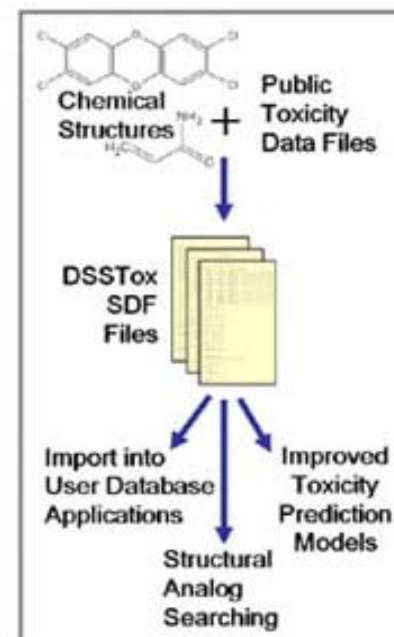
**Site Map**

**Glossary of Terms**

**Help**

**The Distributed Structure-Searchable Toxicity (DSSTox) Database Network provides a community forum for publishing standard format, structure-annotated chemical toxicity data files for open public access. Project goals are to:**

● Encourage use of DSSTox Standard Chemical Structure Fields and SDF standard format files for publishing chemical toxicity databases;

● Coordinate with outside public efforts to encourage chemical structure annotation, data standardization, and open public access to toxicity databases;

● Involve the user community in the effort to migrate more public toxicity data into the DSSTox standardized format for publishing;

● Provide full, open access to toxicity data files for structure-analog searching and for facilitating development of improved models for predicting toxicity based on chemical structure.

Chemical Structures + Public Toxicity Data Files

DSSTox SDF Files

Import into User Database Applications

Improved Toxicity Prediction Models

Structural Analog Searching

DSSTox Graphic Flowchart

**Distributed:** Decentralized set of standardized, field-delimited databases, each separately authored and maintained, that are able to accommodate diverse chemical toxicity data content;
**Structure-Searchable**: Standard format (SDF) structure-data files that can be readily imported into available Chemical Relational Databases and structure-searched;
**Tox:** Toxicity data as it exists in widely disparate forms in current public databases, spanning diverse toxicity endpoints, test systems, levels of biological content, degrees of summarization, and information content.
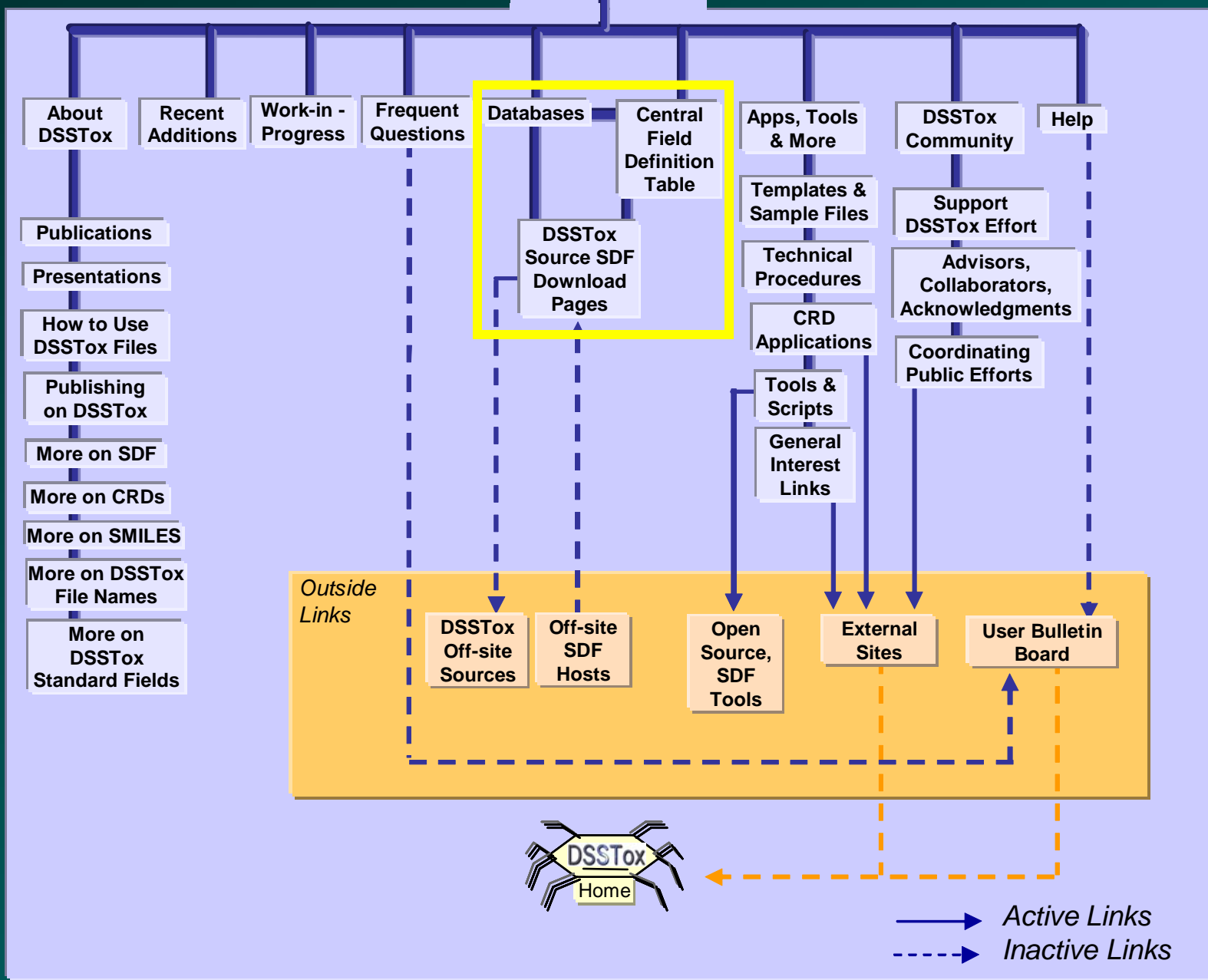
**Current Database Offerings:**
CPDBRM,CPDBHA,CPDBDG,CPDBPR*
DBPCAN
EPAFHM
NCTRER

* CPDBPR last updated 29Mar04

# DSSTox Site Map:

**DSSTox Home**

- **About DSSTox**
  - Publications
  - Presentations
  - How to Use DSSTox Files
  - Publishing on DSSTox
  - More on SDF
  - More on CRDs
  - More on SMILES
  - More on DSSTox File Names
  - More on DSSTox Standard Fields
- **Recent Additions**
- **Work-in - Progress**
- **Frequent Questions**
- **Databases**
  - **Central Field Definition Table**
  - **DSSTox Source SDF Download Pages**
- **Apps, Tools & More**
  - Templates & Sample Files
  - Technical Procedures
  - CRD Applications
  - Tools & Scripts
  - General Interest Links
- **DSSTox Community**
  - Support DSSTox Effort
  - Advisors, Collaborators, Acknowledgments
  - Coordinating Public Efforts
- **Help**

**Outside Links**

- DSSTox Off-site Sources
- Off-site SDF Hosts
- Open Source, SDF Tools
- External Sites
- User Bulletin Board

**DSSTox Home**

→ Active Links

⇢ Inactive Links

## U.S. Environmental Protection Agency

# Distributed Structure-Searchable Toxicity (DSSTox) Public Database Network

Recent Additions | Contact Us | Print Version    Search: [        ] **GO**

EPA Home > Research & Development > National Health and Environmental Effects Research Laboratory > Distributed Structure-Searchable Toxicity (DSSTox) Public Database Network

**About DSSTox**

**Work in Progress**

**Frequent Questions**

**Databases**

**Central Field Definition Table**

**Apps, Tools & More**

**DSSTox Community**

**Site Map**

**Glossary of Terms**

**Help**

**Databases**

- **CPDBRM, CPDBHA, CPDBDO, CPDBPR: Carcinogenic Potency Database Summary Tables for Rat&Mouse, Hamster, Dog, and Non-human Primates**
Tumor target site incidence and TD50 potencies for 1354 chemical substances tested in rats and mouse, 80 chemical substances tested in hamsters, 5 chemicals tested in dogs, and 27 chemical substances tested in non-human primates; data reviewed and compiled from literature and NTP studies.
*(SDF last updated 15Oct03)*

- **DBPCAN: Water Disinfection By-Products Database with Carcinogenicity Estimates**
Carcinogenicity estimates (high, moderate, low concern) by EPA experts using a mechanism-based analog SAR approach on a set of 209 water disinfection by-products, mostly small halogenated organics.
*(SDF last updated 12Sep03)*

- **EPAFHM: EPA Fathead Minnow Aquatic Toxicity Database**
Acute toxicities of 617 chemicals tested in common assay, with mode-of-action assessments.and confirmatory measures.
*(SDF last updated 15Oct03)*

- **NCTRER: FDA's National Center for Toxicological Research - Estrogen Receptor Binding Database**
Estrogen receptor relative binding affinities tested in a common in vitro assay for 232 chemicals, listed with chemical class-based structure activity features.
*(SDF last updated 7Nov03)*

# DSSTox Database Design:

# DSSTox Toxicity Database Standards:

- Data file format (SDF)

- File naming convention

- Chemical structure information fields

- Documentation requirements

- Publishing requirements

# DSSTox Standard Chemical Fields:

- Structure          *2D chemical structure*
- StructureShown          *Description of displayed 2D structure*
  - *tested form, simplified to parent, predicted form, active ingredient of formulation, general form*
- Formula          *Empirical molecular formula*
- MolWeight          *Molecular weight in atomic units*
- CAS          *Chem Abstracts Service No. for StructureShown*
- SMILES          *Linear text notation for 2D StructureShown*
- DSSTox_ID          *Counter allows unique identification of record*
- DSSTox_FileName          *Name of file included in each record*
- ChemName          *Chemical name from original data base*
- SubstanceType          *Broad substance classification*
  - *defined organic, inorganic, organometallic, polymer, mixture or unknown*
- TestedForm          *Tested form of chemical*
  - *parent, salt, complex, unknown or multiple forms*
- AddToParent          *Salt counterions or complexed moieties*
- CAS_TestedForm          *CAS No. for tested form of chemical*
- SMILES_TestedForm          *SMILES code for the tested form of the chemical*
- ChemNote          *Additional qualifier info for chemical fields*
  - *defined mixture characteristics, uncertainty in structure or CAS, stereochem, replicate, etc.*
- ChemCount          *Counter for structure or CAS duplications in db*

Structure

Structure1

**StructureShown**

Tested form
Simplified to parent
General form
Active ingredient of formulation
Monomer of polymeric form

*Linked-content fields*

CAS
SMILES
Formula
MolWt

# Substance Type

NO Structure OR structures unassociated species

**mixture or unknown**

* No structure
* No SMILES
* No CAS Parent
* ChwmNote field populated

NO

2 or more associated species

NO

YES

Single defined structure

NO

**complex**

*(Tested Form)*

YES

**inorganic**

NO

Contains Carbon

YES

**defined organic**

*(Only class in DOP file)*

NO

Contains metal or metaloid

YES

**organometallic**

*complex*

**organometallic**

*parent complex salt*

**defined organic**

*parent complex salt*

**inorganic**

| H | | | | | | | | | | | | | | | | | He |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Li | Be | | | | | | | | | | | B | C | N | O | F | Ne |
| Na | Mg | | | | | | | | | | | Al | Si | P | S | Cl | Ar |
| K | Ca | Sc | Ti | V | Cr | Mn | Fe | Co | Ni | Cu | Zn | Ga | Ge | As | Se | Br | Kr |
| Rb | Sr | Y | Zr | Nb | Mo | Tc | Ru | Rh | Pd | Ag | Cd | In | Sn | Sb | Te | I | Xe |
| Cs | Ba | | Hf | Ta | W | Re | Os | Ir | Pt | Au | Hg | Tl | Pb | Bi | Po | At | Rn |
| Fr | Ra | | Rf | Db | Sg | Bh | Hs | Mt | Uun | Uuu | Uub | | | | | | |

salts

metals

| La | Ce | Pr | Nd | Pm | Sm | Eu | Gd | Tb | Dy | Ho | Er | Tm | Yb | Lu |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ac | Th | Pa | U | Np | Pu | Am | Cm | Bk | Cf | Es | Fm | Md | No | Lr |

# Integrating Diverse Databases from a Chemical Structure Perspective:

CPDB        DBPCAN        EPAFHM        NCTRER        ….

## Standard Chemical Fields

**SAL CPDB**

**TD50 Rat**

**TD50Mouse**

**Target Sites Rat Male**

**Target Sites Rat Female**

**Target Sites Mouse Male**
**….**
**Other Species**

**ChemClass DBP**

**Concern Level**

**Rationale**

**Rational Source**

**Analog ChemName**

**AnalogCAS**

**AnalogSMILES**

**ChemClass FHM**

**MOA**

**MOACONF**

**CLOGP**

**LC50**

**LC50NOTE**

**LC50RATIO**

**MIXMOA**

**TOXINDEX**

**FATS**

**BEHAVIOR**

**NCTRlogRBA**

**ER RBA**

**ChemClass ERB**

**Activity Group ERB**

**Rationale ChemClass ERB**

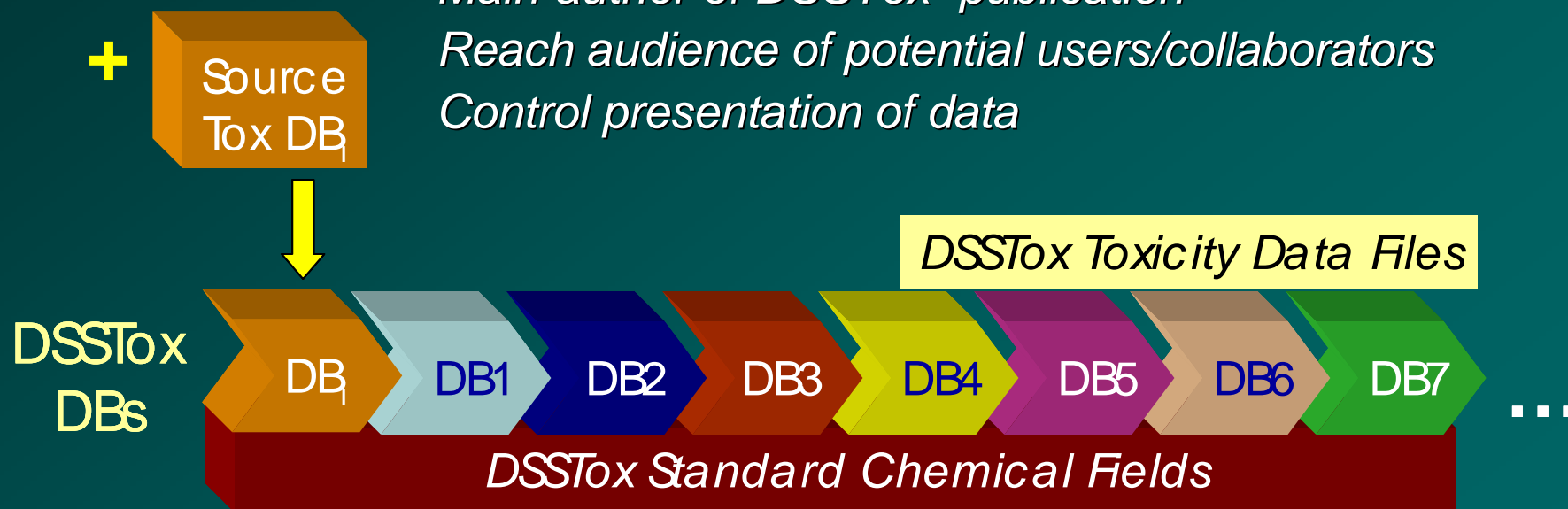**MeanChem Class ERB RBA**

**LogP**

**F1, F2, …F6**

# DSSTox Database Network:

➢ How can DSSTox files be used?

➢ What's next?

# Begin to incorporate standard tox fields (ToXML)

## CPDB　　　DBPCAN　　　EPAFHM　　　NCTRER　　.....

**Standard Chemical Fields**

*Standard Tox Fields:*　species, sex, strain, assay, dose

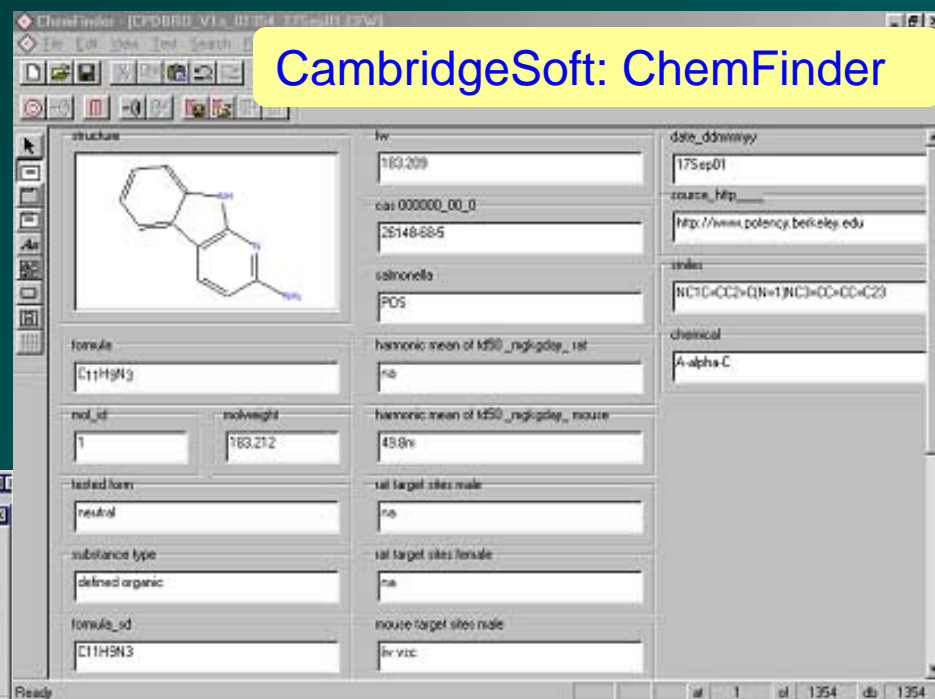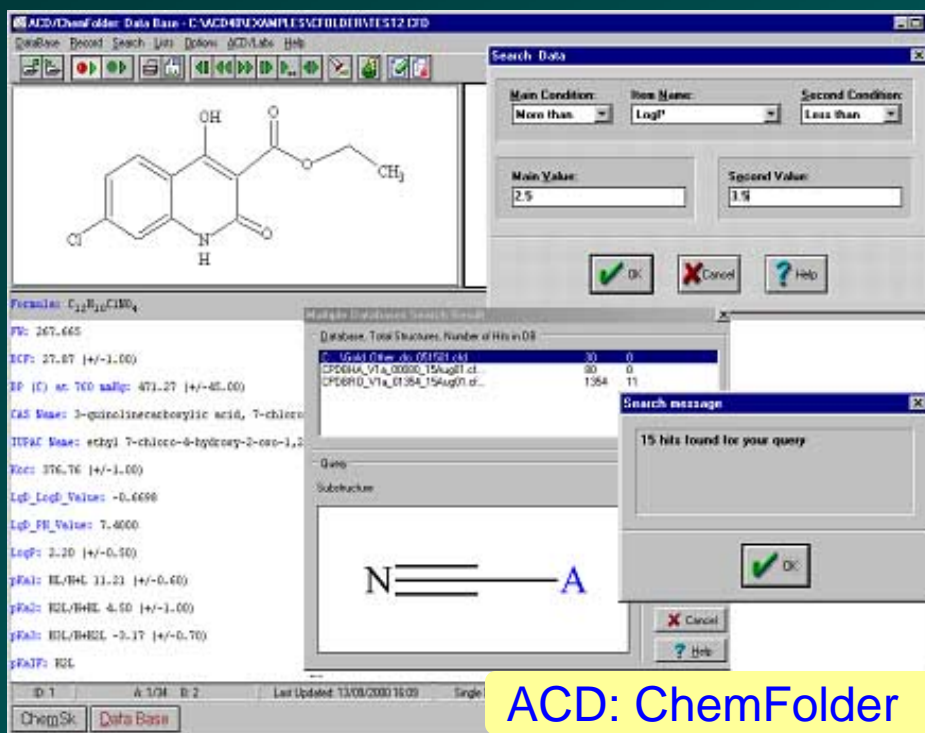| | | | |
|---|---|---|---|
| **SAL CPDB** | | **ChemClass FHM** | |
| | **ChemClass DBP** | **MOA** | **NCTRlogRBA** |
| **TD50 Rat** | | **MOACONF** | |
| | **Concern Level** | **CLOGP** | **ER RBA** |
| **TD50Mouse** | | **LC50** | |
| | **Rationale** | **LC50NOTE** | **ChemClass ERB** |
| **Target Sites Rat Male** | **Rational Source** | **LC50RATIO** | **Activity Group ERB** |
| **Target Sites Rat Female** | **Analog ChemName** | **MIXMOA** | **Rationale ChemClass ERB** |
| | **AnalogCAS** | **TOXINDEX** | |
| **Target Sites Mouse Male** | **AnalogSMILES** | **FATS** | **MeanChem Class ERB RBA** |
| **....** | | **BEHAVIOR** | **LogP** |
| **Other Species** | | | **F1, F2, …F6** |

# Migrate More Public Toxicity Data into DSSTox Standard Format: Phase II, III, …

- FDA Max Recommended Therapeutic Dose - Pharmaceuticals

- NCTR Androgen, Thyroid, and Endocrine Disruption Databases

- NTP Rodent carcinogenicity bioassays, subchronic bioassays, developmental, repro, immuno, etc.

- ICVAM databases on LD50, skin sensitization, local lymph node assay, skin corrosivity, endocrine disruption, etc

- EPA pesticide ecotoxicity database

- Developmental toxicity database (literature – FDA,TOPKAT)

- UniLever Skin Sensitization database

- EPA's High Production Volume (HPV) chemical data

- EPA's Integrated Risk Information System (IRIS)

# Chemical Relational Databases: *Exploration across toxicological domains and structural/biological axes*

Accord
Oracle
ISIS
ChemFolder
ChemFinder
LeadScope
BioRad

CambridgeSoft: ChemFinder

ACD: ChemFolder

Freeware SDF Viewer Application

Non-EPA on-line structure searching

EPA on-line structure searching

# Coordinating Public Efforts:

DSSTox **COMMUNITY**

- ACD/Labs (Advanced Chemistry Development) ChemFolder Public Databases
- Cambridge-Soft's ChemFinder.Com Chemical Search Website
- FDA (Food & Drug Administration) Center for Drug Evaluation & Research
- ILSI (International Life Sciences Institute) SAR Toxicity Database Project, in collaboration with LHASA, Lmt.LIST (LeadScope In Silico Tox) Focus Group
- LIST (LeadScope In Silico Tox) Focus Group
- MGED: MIAMI-Tox
- NCI (National Cancer Institute) Public Data Outreach – Structure Web Browser
- NIEHS's National Center for Toxicogenomics
- NLM (National Library of Medicine) TOXNET
- NTP (National Toxicology Program) On-line Public Databases
- SRC (Syracuse Research Corporation) PBT-Profiler and Analog Search Tools

**NCI Structure-Browser**

**Advanced Chemistry Development**

**ISSCentral – Freeware CRD**

**EPA/Syracuse Res Corp**

# Collaboration with NIEHS/National Center for Toxicogenomics: *Chemical Effects in Biological System Knowledge-Base (CEBS)*

Gene expression

Gene Pathways

Gene Function

Proteomics

Historical Toxicity data



## CEBS Vision - Bioinformatics to Knowledge

**Associated Data**

DATABASE OF ARRAY, PROTEOMICS AND TOXICOLOGY DATA ON CHEMICALS, DRUGS AND STRESSORS

DATABASE ON GENES AND GENE GROUPS RELEVANT TO ENVIRONMENTAL DISEASE

DATABASE OF SNPs AND MUTANTS RELEVANT TO ENVIRONMENTAL DISEASE

RETRIEVE

COMPENDIA OF FUNCTIONAL GENE GROUPS WITH ASSOCIATED PATHWAYS AND NETWORKS

DICTIONARIES AND METADATA

STORE AND CONVERT

**External Links**

NTP AND OTHER TOX. DATABASES

LINKS

PATHWAYS

FUNCTION

PROTEIN DBs

NLM/NCBI

GENOMIC RESOURCES

GENE/PROT DESCRIPTIONS

**Query**

COMPOUND/CLASS/STRUCTURE

EFFECTS

GENES/PROTEIN FUNCTIONAL GROUPS

# Collaboration with NIEHS/National Center for Toxicogenomics: *Chemical Effects in Biological System Knowledge-Base (CEBS)*
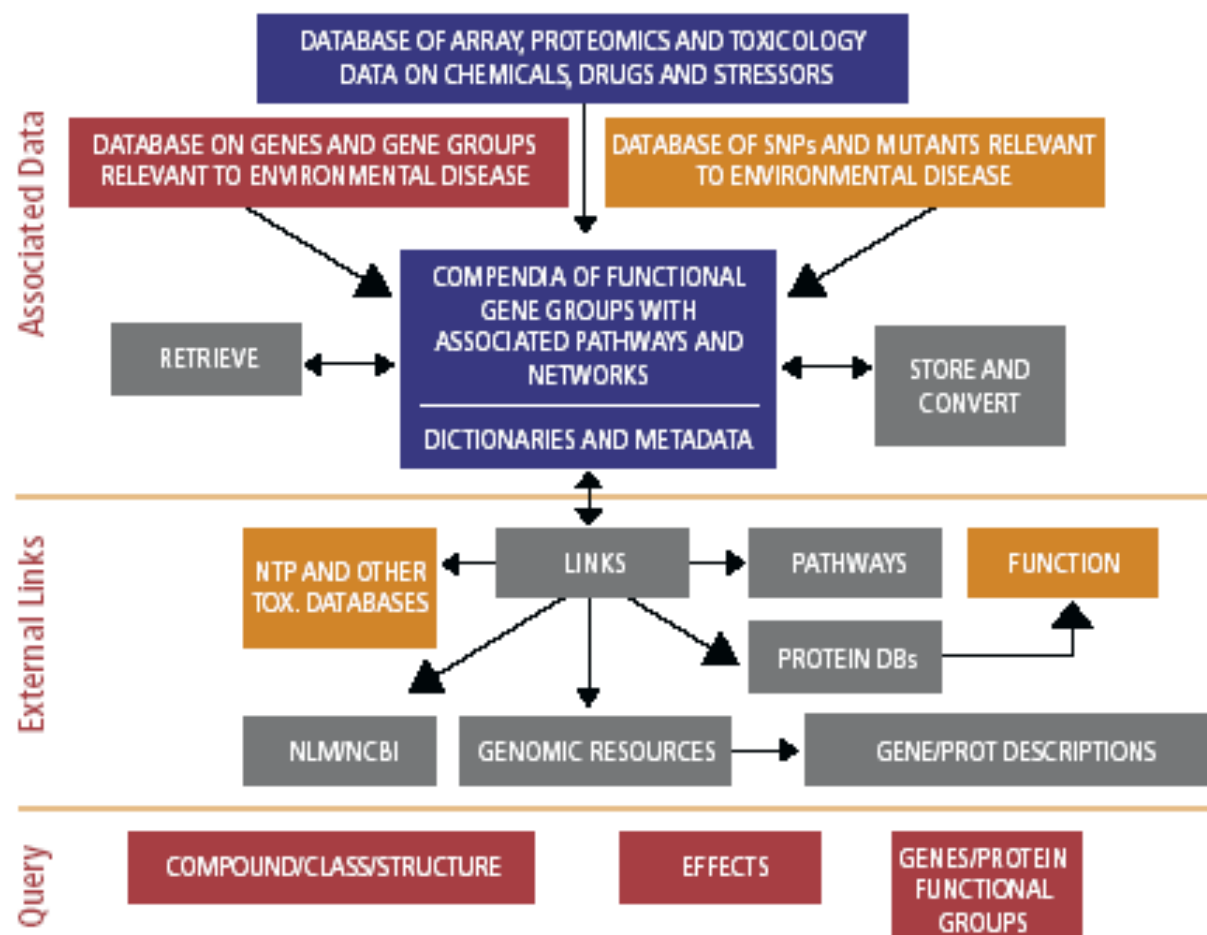
**DSSTox +**

Gene expression

Gene Pathways

Gene Function

Proteomics

*DSSTox Toxicity Data Files*

Historical Toxicity data

DB1  DB2  DB3  DB4  DB5  DB6  DB7

*DSSTox Standard Chemical Fields*



## CEBS Vision - Bioinformatics to Knowledge

DATABASE OF ARRAY, PROTEOMICS AND TOXICOLOGY DATA ON CHEMICALS, DRUGS AND STRESSORS

DATABASE ON GENES AND GENE GROUPS RELEVANT TO ENVIRONMENTAL DISEASE

DATABASE OF SNPs AND MUTANTS RELEVANT TO ENVIRONMENTAL DISEASE

**Associated Data**

RETRIEVE

COMPENDIA OF FUNCTIONAL GENE GROUPS WITH ASSOCIATED PATHWAYS AND NETWORKS

DICTIONARIES AND METADATA

STORE AND CONVERT

**External Links**

NTP AND OTHER TOX. DATABASES

LINKS

PATHWAYS

FUNCTION

PROTEIN DBs

NLM/NCBI

GENOMIC RESOURCES

GENE/PROT DESCRIPTIONS

**Query**

COMPOUND/CLASS/STRUCTURE

EFFECTS

GENES/PROTEIN FUNCTIONAL GROUPS

# Chemo-bioinformatics:
## Expanded Relational Search Domains

Historical Toxicity Data

Species
Activity
Tissue
Dose

Genes
Gene families
Proteins
Profiles

Functional groups

Phys-chem properties

Structural analogs

Size
Shape

Chemical Structures

Genomics Proteomics

# DSSTox Collaborators/Advisors/Acknowledgements:

- Cancer Potency Data Base - rodent carcinogenicity……… Lois Swirsky Gold, Thomas Slone
- EPA - Ecotoxicity (fathead minnow, Teratox)……………… Chris Russom
- EPA/OPPT,OW – DBP cancer assessment……………… Yin-tak Woo, Mary Manibusan
- FDA/NCTR - Estrogen receptor binding data base………. Weida Tong, Hong Fang
- FDA/CDER – Maximum recommended dose drugs …….. Dan Benz, Ed Matthews, Joe Contrera
- EPA/OPP, Ecotox – pesticides ……………………………… Brian Montague, Pauline Wagner
- NTP Gene-tox data; IRIS ……………………………………. Errol Zeiger, Zeiger Consulting
- Nat Ctr Toxicogenomics: CEBS ………………………….. Mike Waters, Ray Tennant
- GlaxoSK - GeneTox/NTP Salmonella database………… Neal Cariello, Vijay Gombar
- NIEHS/NTP Rodent carcinogenicity, etc ………………… Skip Eastin, Doug Bristol
- Developmental toxicity …………………………………………. Vijay Gombar (GlaxoSK), Orest Macina
- ICVAM Toxicity databases ……………………………………. Ray Tice, Marc Jackson, ILS
- Unilever Skin Sensitization Database……………………… Martin Barrett, Marlin Consulting
- ZEBET Acute Toxicity Database ……………………………. Julie Penzotti, Rational Discovery
- Tulane/Xavier Univ – Endocrine Disruption ……………… Tom Wiese
- NCI – SDF tools, CACTVS structure browser……………. Marc Nicklaus
- LeadScope – SDF/XML converter, FDA carcinogens …….. Chihae Yang
- SDF Viewer application ………………………………………. Thomas Harrocks, Intuitive Software Solutions
- ACD – ChemFolder, WebBase ……………………………….. Antony Williams, M Hachey, G Shear
- CambridgeSoft – ChemFinder application ……………….. Rich Talbot
- EPA scientific advisors ……………………………………….. Stephen Nesnow, Adam Swank
- EPA/CSC – web development, IT …………………………… Brian Garges/ D Kanipe, D Marshall

# DSSTox Development Team:

- ClarLynda Williams – Coordinator, Lead Project Technician
  - *NCCU/EPA COOP Student Trainee 12/00-12/02; EPA 12/02-8/03*
- Jamie Burch
  - *NCCU/EPA COOP Student Trainee 11/02-2/04*
- Brian Rogers
  - *NCCU/EPA COOP Student Trainee 1/04-present*
- Audrey Evans
  - *ECO/EPA Summer Student Trainee, Summer '03*
- Todd Stewart
  - *EPA-UNC Student COOP, Spring/Summer 02*
- Nina Fields
  - *Shaw Univ. High School Minority Mentoring Program, Summer '02*
- James Beidler
  - *EPA Summer Student Employee, Warren-Wilson College, Summer '01*
- Daniel Ohuoba
  - *Shaw Univ. High School Minority Mentoring Program, Summer '01*
- Adam Swank
  - *ECD, EPA*